

Accurate assembly of the olive baboon (*Papio anubis*) genome using long-read and Hi-C data

--Manuscript Draft--

Manuscript Number:	GIGA-D-20-00013	
Full Title:	Accurate assembly of the olive baboon (<i>Papio anubis</i>) genome using long-read and Hi-C data	
Article Type:	Data Note	
Funding Information:	National Institutes of Health (US) (R24 OD017859)	Dr. Laura A. Cox Dr. Jeffrey D. Wall
	National Institutes of Health (US) (R01 GM094402)	Dr. Yun S. Song
	National Institutes of Health (US) (R01 HG005946)	Dr. Pui-Yan Kwok
	David and Lucile Packard Foundation (Packard Fellowship for Science and Engineering)	Dr. Yun S. Song
Abstract:	<p>Background</p> <p>Besides macaques, baboons are the most commonly used nonhuman primate in biomedical research. Despite this importance, the genomic resources for baboons are quite limited. In particular, the current baboon reference genome Panu_3.0 is a highly fragmented, reference-guided (i.e., not fully <i>de novo</i>) assembly, and its poor quality inhibits our ability to conduct downstream genomic analyses.</p> <p>Findings</p> <p>Here we present a truly <i>de novo</i> genome assembly of the olive baboon (<i>Papio anubis</i>) that uses data from several recently developed single-molecule technologies. Our assembly, Panubis1.0, has an N50 contig size of ~1.46 Mb (as opposed to 139 Kb for Panu_3.0), has single scaffolds that span each of the 20 autosomes and the X chromosome, and is freely available for scientific use from NCBI.</p> <p>Conclusions</p> <p>We present multiple lines of evidence (including Bionano Genomics data, linkage information, and patterns of linkage disequilibrium) suggesting that the Panubis1.0 assembly corrects large assembly errors in Panu_3.0. This in turn has led to an improved baboon annotation, making Panubis1.0 much more useful for future genomic studies.</p>	
Corresponding Author:	Jeffrey D. Wall UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Sanjit Singh Batra	
First Author Secondary Information:		
Order of Authors:	Sanjit Singh Batra	
	Michal Levy-Sakin	
	Jacqueline Robinson	

	Joseph Guillory
	Steffen Durinck
	Tauras P. Vilgalys
	Pui-Yan Kwok
	Laura A. Cox
	Somasekar Seshagiri
	Yun S. Song
	Jeffrey D. Wall
Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
Availability of data and materials	Yes

All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

Accurate assembly of the olive baboon (*Papio anubis*) genome using long-read and Hi-C data

Sanjit Singh Batra¹, Michal Levy-Sakin², Jacqueline Robinson³, Joseph Guillory⁴, Steffen Durinck^{4,5}, Tauras P. Vilgalys⁶, Pui-Yan Kwok^{2,3}, Laura A. Cox^{7,8}, Somasekar Seshagiri⁴, Yun S. Song^{1,9,10} and Jeffrey D. Wall^{3,*}

¹Computer Science Division, University of California, Berkeley, CA 94720;

²Cardiovascular Research Institute, University of California San Francisco, San Francisco, CA 94143;

³Institute for Human Genetics, University of California San Francisco, San Francisco, CA 94143;

⁴Department of Molecular Biology, Genentech, Inc., 1 DNA Way, South San Francisco, CA 94080;

⁵Bioinformatics and Computational Biology Department, Genentech, Inc., 1 DNA Way, South San Francisco, CA 94080;

⁶Department of Evolutionary Anthropology, Duke University, Durham, NC 27705

⁷Center for Precision Medicine, Department of Internal Medicine, Section of Molecular Medicine, Wake Forest School of Medicine, Winston-Salem, NC 27101

⁸Southwest National Primate Research Center, Texas Biomedical Research Institute, San Antonio, TX 78245

⁹Department of Statistics, University of California, Berkeley, CA 94720;

¹⁰Chan Zuckerberg Biohub, San Francisco, CA 94158

*Corresponding Author: Jeff.Wall@ucsf.edu

1 **ABSTRACT**

2

3 **Background**

4

5 Besides macaques, baboons are the most commonly used nonhuman primate in
6 biomedical research. Despite this importance, the genomic resources for baboons are
7 quite limited. In particular, the current baboon reference genome Panu_3.0 is a highly
8 fragmented, reference-guided (i.e., not fully *de novo*) assembly, and its poor quality
9 inhibits our ability to conduct downstream genomic analyses.

10

11 **Findings**

12

13 Here we present a truly *de novo* genome assembly of the olive baboon (*Papio anubis*)
14 that uses data from several recently developed single-molecule technologies. Our
15 assembly, Panubis1.0, has an N50 contig size of ~1.46 Mb (as opposed to 139 Kb for
16 Panu_3.0), has single scaffolds that span each of the 20 autosomes and the X
17 chromosome, and is freely available for scientific use from NCBI.

18

19 **Conclusions**

20

21 We present multiple lines of evidence (including Bionano Genomics data, linkage
22 information, and patterns of linkage disequilibrium) suggesting that the Panubis1.0
23 assembly corrects large assembly errors in Panu_3.0. This in turn has led to an
24 improved baboon annotation, making Panubis1.0 much more useful for future genomic
25 studies.

26 Data Description

28 Introduction

29
30 Baboons are ground-living monkeys native to Africa and the Arabian Peninsula. Due to
31 their relatively large size, abundance and omnivorous diet, baboons have increasingly
32 become a major biomedical model system (reviewed in [1]). Baboon research has been
33 facilitated by the creation (in 1960) and maintenance of a large, pedigreed, well-
34 phenotyped baboon colony at the Southwest National Primate Research Center
35 (SNPRC) and an ability to control the environment of subjects in ways that are obviously
36 not possible in human biomedical studies. For example, baboons have been used to
37 study the effect of diet on cholesterol and triglyceride levels in controlled experiments
38 where all food consumption is completely controlled [2] [3] [4]. In recent years, linkage
39 studies in baboons have helped identify genetic regions affecting a wide range of
40 phenotypes, such as cholesterol levels [5] [6], estrogen levels [7], craniofacial
41 measurements [8], bone density [9] [10] and lipoprotein metabolism [11]. In addition,
42 studies have also documented that the genetic architecture of complex traits in baboons
43 can be directly informative about analogous traits in humans (e.g., [10] [12]).

44 The success of these and other studies have been mediated in part by recent
45 advances in molecular genetics technologies. In particular, the ability to cheaply
46 genotype and/or sequence samples of interest has led to a revolution in genetic studies
47 of the associations between genotype and phenotype. While human genetic studies
48 now routinely include the analyses of whole-genome sequence data from many
49 thousands of samples (e.g., [13] [14] [15] [16] [17]), comparable studies in other
50 organisms have lagged far behind. Part of the reason for this is the lack of genetic
51 resources in non-human species. Large, international projects such as the Human
52 Genome Project [18] [19], International HapMap Project [20] [21] [22] and the 1000
53 Genomes Project [23] [24] [25] have provided baseline information on sequences and
54 genetic variation, while complementary efforts to quantify recombination rates have led
55 to detailed genetic maps (PMIDs: 8600387, 12053178, 16224025, 20981099).
56 Subsequent human genetic studies have utilized all of this background information.

57 The first published baboon genome assembly was from a yellow baboon [26]. This
58 assembly used a combination of Illumina paired-end and Illumina mate-pair sequence
59 data (with mean library insert sizes ranging from 175 bp to 14 Kbp) to produce a highly
60 fragmented assembly with contig N50 of 29 Kbp and scaffold N50 of 887 Kbp. The
61 public olive baboon assembly, Panu_3.0, suffers from the same problem of having small
62 contigs and scaffolds (contig N50 of 139 Kbp and *de novo* scaffold N50 of 586 Kbp)
63 [27]. The authors of the public olive baboon assembly chose to distribute a reference-
64 guided assembly with scaffolds mapped onto rhesus (*Macaca mulatta*) chromosomes.
65 As a consequence, many syntenic differences between rhesus and baboon will result in
66 large-scale assembly errors in Panu_3.0. One additional drawback of this baboon
67 genome assembly was its informal embargo from 2008 to 2019 under the guidelines of
68 the Fort Lauderdale agreement. Hence, its influence on scientific research has been
69 negligible.

70 In this project, we provide a high-quality, *de novo* genome assembly for olive
71 baboon (*Papio anubis*), which we call Panubis1.0, with the hope that this resource will

72 enable future high-resolution genotype-phenotype studies. Unlike previous baboon
73 genome assembly efforts, we use a combination of three recently developed
74 technologies (from 10x Genomics linked-reads, Oxford Nanopore long reads, and Hi-C)
75 to increase the long-range contiguity of our assembly. These newly developed
76 technologies enable us to generate assemblies where the autosomes (and the X
77 chromosome) are each spanned by a single scaffold at a cost that is orders of
78 magnitude cheaper than the Panu_3.0 assembly. We also verify that most of the large-
79 scale syntenic differences between our Panubis1.0 and Panu_3.0 are due to errors in
80 Panu_3.0 rather than Panubis1.0. Our assembly is available for scientific use without
81 any restrictions.

82

83 **Genome Sequencing**

84

85 Index animal: We used individual number 15944 (currently deceased) from the SNPRC
86 pedigreed baboon colony for all of the sequencing and genome assembly work
87 associated with this project.

88

89 10x Genomics sequencing: High molecular weight genomic DNA extraction, sample
90 indexing, and generation of partition barcoded libraries were performed according to the
91 10x Genomics (Pleasanton, CA, USA) Chromium Genome User Guide and as
92 published previously ([28]). An average depth of ~60X was produced and analyzed for
93 this project.

94

95 Oxford Nanopore sequencing: Libraries for the Oxford Nanopore sequencing were
96 constructed as described previously ([45]) using DNA derived from whole blood. The
97 sequencing was conducted at Genentech, Inc. (South San Francisco, CA, USA); we
98 analyzed data with an average depth of ~15X for this project.

99

100 Bionano optical maps: High-molecular-weight DNA was extracted, nicked, and labeled
101 using the enzyme Nt.BspQI (New England Biolabs (NEB), Ipswich, MA, USA), and
102 imaged using the Bionano Genomics Irys system (San Diego, CA, USA) to generate
103 single-molecule maps for assessing breaks in synteny between Panu_3.0 and
104 Panubis1.0.

105

106 Hi-C sequencing: High molecular weight DNA from Jenny Tung (Duke University) was
107 sent to Phase Genomics. ~15X Hi-C data was obtained using previously described
108 techniques [46].

109

110 **Linkage disequilibrium analyses**

111

112 We estimated the scaled recombination rate ρ ($= 4Nr$ where N is the effective
113 population size and r is the recombination rate per generation) using LDhelmet [47] from
114 24 unrelated olive baboons [48]. We then identified potential breaks in synteny as
115 regions with total $\rho > 500$ and $\rho / \text{bp} > 0.2$. We considered there to be evidence of a
116 synteny break if one of these regions was within 50 Kb of a potential breakpoint (as

117 identified in Panu_3.0 vs. Panubis1.0 comparisons). The false discovery rate, based on
118 randomly permuting the locations of LD-based potential breaks in synteny, is ~4%.

119
120 To calculate recombination rates, we used a variant call set mapped onto the old
121 assembly Panu_2.0, as described in [48]. For the potential breaks in synteny identified
122 above, we used liftover to convert the breakpoints into Panu_3.0 coordinates and
123 verified that Panu_2.0 and Panu_3.0 were syntenic with each other across the
124 breakpoints.

125
126 Finally, due to the inherent noise in linkage-disequilibrium based estimates of ρ , the lack
127 of evidence for a synteny break in Panu_3.0 is not positive evidence that the Panu_3.0
128 assembly is correct.

129 130 **Inference of crossovers in a baboon pedigree**

131
132 We utilized a previously described vcf file for the baboons shown in Figure 6 which was
133 mapped using Panu_2.0 coordinates and lifted over to Panu_3.0 coordinates. We
134 filtered for high quality genotype calls, considered only biallelic SNPs, and required a
135 depth ≥ 15 , QUAL > 50 and genotype quality (GQ) ≥ 40 in order to make a genotype
136 call. We further required an allelic balance (AB) of > 0.3 for heterozygote calls and AB
137 < 0.07 for homozygote calls, and excluded all repetitive regions as described in [48].

138
139 We focused our analyses on those SNPs that were most informative about recent
140 crossover events. For example, to detect paternal crossovers, we restricted our
141 analyses to SNPs where the father (i.e., 10173 in Figure 6) was heterozygous, both
142 mothers (9841 and 12242 in Figure 6) were homozygous, and all 9 offspring had
143 genotype calls. (For maternal crossovers, we required 10173 to be homozygous and
144 both 9841 and 12242 to be heterozygous.) For these sites, it is straightforward to infer
145 which allele (coded as 0 for reference allele and 1 for alternative allele) was passed on
146 from 10173 to his offspring. While the haplotypic phase of 10173 is unknown, we can
147 infer crossover events based on the minimum number of crossovers needed to be
148 consistent with the observed patterns of inheritance in the offspring of 10173 ([49]). For
149 example, Figure 5c shows that the inheritance pattern near position 29.38 requires at
150 least 3 crossovers (e.g., in individuals 17199, 18385 and 19348).

151
152 For each potential error in the Panu_3.0 assembly, we converted the breakpoint
153 location into Panu_2.0 coordinates and verified synteny between Panu_2.0 and
154 Panu_3.0 across the breakpoint region. We then determined whether there were an
155 abnormally large number of crossovers inferred right at the breakpoint. Specifically, if
156 we inferred at least 3 crossover events (out of 18 total meioses, 9 paternal and 9
157 maternal), then we considered this as evidence that the Panu_3.0 assembly is incorrect,
158 as in Figure 5c (cf. 'Linkage Support' column in Table 3). Note that the converse isn't
159 true: fewer than 3 inferred crossover events is not evidence that the Panu_3.0 assembly
160 is correct at a particular location.

161 162 **Genome Assembly**

163
164 The main strength of our approach is in combining data from multiple platforms (10x
165 Genomics linked-reads, Oxford Nanopore long-reads, Illumina paired-end short-reads,
166 and Hi-C), which have complementary advantages (Figure 1). Figure 1 describes our
167 assembly strategy. We began by assembling 10x Genomics reads generated with their
168 Chromium system (average depth ~60x) using the SUPERNOVA assembler (version
169 1.1) [28], which yielded an assembly with a contig N50 of ~84 kb and a scaffold N50 of
170 ~15.7 Mb (Table 1). The gap lengths between the contigs in a scaffold obtained by
171 assembling 10x linked-reads are arbitrary [29]. Hence, in order to leverage the Oxford
172 Nanopore long-reads for gap-closing, we split the 10X scaffolds at every stretch of non-
173 zero N's to obtain a collection of contigs.

174
175 We scaffolded the resulting contigs with Oxford Nanopore long-reads (average depth
176 ~15X) using the LR_Scaf [30] scaffolding method. This resulted in an assembly with a
177 contig N50 of ~134 kb and a scaffold N50 of ~1.69 Mb (Table 1). These resulting
178 scaffolds are more amenable to gap-closing, because the gap lengths (number of Ns
179 between two consecutive contigs) are estimated by long-reads that span each gap and
180 align to the flanking regions of that gap.

181
182 Upon performing gap-closing with the same set of Oxford Nanopore long-reads using
183 LR_Gapcloser [31], we obtained an assembly with a contig N50 of ~1.47 Mb and a
184 scaffold N50 of ~1.69 Mb (Table 1). Note that this increase in contig N50 of ~84Kb from
185 the 10x Genomics linked-read assembly, to a contig N50 of ~1.47 Mb, would not have
186 been possible if we had simply performed gap-closing with the Oxford Nanopore long
187 reads directly on the 10x-based assembly without first splitting it into its constituent
188 contigs. Finally, we polished the resulting assembly by aligning Illumina paired-end
189 reads (average depth ~60X in PE150 reads) using Pilon [32].

190
191 In order to scaffold the resulting assembly with Hi-C data, we first set aside scaffolds
192 shorter than 50 kb, which comprised only ~1.8% of the total sequence base pairs. This
193 was done because Hi-C based scaffolding is more reliable for longer scaffolds, since
194 there are more Hi-C reads aligning to longer scaffolds. We then ordered and oriented
195 the remaining scaffolds using the 3D *de novo* assembly (3d-dna) pipeline [33] using
196 ~15X Hi-C data generated by Phase Genomics [34]. Finally, we manually corrected
197 misassemblies in the resulting Hi-C based assembly by visualizing the Hi-C reads
198 aligned to the assembly, using Juicebox Assembly Tools [35], following the strategy
199 described in [36]. Figure 2 shows Hi-C reads aligned to the resulting assembly with the
200 blue squares on the diagonal representing chromosomes.

201
202 The resulting *Papio anubis* genome assembly, which we name Panubis1.0, contains
203 ~2.87 Gb of sequenced base pairs (non-N base pairs) and 2.3 Mb (<0.1%) of gaps
204 (N's). Single scaffolds spanning the 20 autosomes and the X chromosome together
205 contain 95.14% (~2.73 Gb) of the sequenced base pairs. We number the autosomes as
206 chr1 to chr20, in decreasing order of the scaffold length, so some chromosome
207 numbers in our convention are different from Panu_3.0's numbering. We note that
208 Panubis1.0 has a contig N50 of 1.46 Mb, which is a greater than ten-fold improvement

209 over the contig N50 (~139 kb) of the Panu_3.0 assembly. As a result, Panubis1.0
210 contains five times fewer scaffolds (12,976 scaffolds with a scaffold N50 of ~140 Mb)
211 compared to the Panu_3.0 assembly (63,235 scaffolds with a scaffold N50 of ~586 Kb);
212 see Table 2 for a further comparison of the two assemblies. Gene completion analysis
213 of the assembly using BUSCO v2 and the odb9 Mammalia ortholog dataset [37]
214 suggests that Panubis1.0 contains 93.00% complete genes, comparable to the
215 Panu_3.0 assembly.

216
217

218 **Y chromosome assembly**

219

220 The Hi-C scaffolding with 3d-dna yielded an ~8 Mb scaffold that putatively represents
221 part of the baboon Y chromosome. Since, rhesus macaque is the phylogenetically
222 closest species to baboons which has a chromosome-scale assembly, we aligned this
223 putative baboon Y chromosome scaffold with the rhesus macaque Y chromosome
224 (Figure 3). We observed a substantial amount of synteny between the putative baboon
225 Y and the rhesus Y, comparable to what is observed between the chimpanzee Y and
226 the human Y chromosomes. (For reference, genetic divergence between baboon and
227 rhesus is similar to human – chimpanzee divergence [38].) The observed breaks in
228 synteny are consistent with the well-documented high rate of chromosomal
229 rearrangements on mammalian Y chromosomes [39].

230

231

232 **Genome Annotation**

233

234 Annotation of the protein and non-protein coding genes was performed by NCBI
235 (O’Leary et al. 2016, *Nucleic Acids Research*, Reference sequence (RefSeq) database
236 at NCBI), based on RNA sequencing of 4 captive baboons at the SNPRC (BioProject
237 PRJNA559725) as well as other publically available baboon expression data.

238 Panubis1.0 contains 21,087 protein-coding genes and 11,295 non-coding genes. This is
239 a slight decrease in the number of protein-coding genes relative to Panu_3.0 (21087 vs
240 21,300) which can be explained by merging genes together (n=252), and an increase in
241 the number of non-coding genes (11295 vs 8433). Panubis1.0 also contains slightly
242 more pseudogenes (6680 vs 5998) and genes with splice variants (14526 vs 13693).

243 Many of these differences may reflect insights gained from an improved assembly
244 leading to an increased ability to map sequencing data; indeed, during genome
245 annotation, 88% of RNA-seq reads mapped to Panubis1.0 while only 80% mapped to
246 Panu3.0.

247

248 Overall, most genes (66%) are highly similar or identical between Panubis1.0 and
249 Panu_3.0. Of remaining genes, 13% of genes contain major changes (e.g. were split,
250 moved, changed gene type, or changed substantially in completeness), 20% are novel
251 in Panubis1.0, and 12% deprecated from Panu_3.0.

252

253 **Comparisons with the publicly available Panu_3.0 assembly**

254

255 We constructed chromosome-scale dotplots to identify large syntenic differences
256 between the Panubis1.0 and Panu_3.0 assemblies (Figure 4) and examined more
257 closely all of the large (>100 Kb) differences between the two (Table 3). We used
258 several orthogonal sources of information to assess whether these differences were
259 errors in our Panubis1.0 assembly or in the Panu_3.0 assembly. These included
260 Bionano Genomics optical maps obtained from the same individual used for generating
261 Panubis1.0, linkage information from a pedigree of baboons that were all sequenced to
262 high coverage, and linkage-disequilibrium information from 24 unrelated olive baboons
263 from the SNPRC pedigreed baboon colony. We manually examined each break in
264 synteny between Panubis1.0 and Panu_3.0 to determine whether these independent
265 sources of evidence supported one assembly over the other (summarized in Table 3).
266 Overall, in 11 out of 12 large syntenic differences between Panubis1.0 and Panu_3.0, at
267 least one of these independent sources provided evidence that the Panubis1.0
268 assembly is correct. These independent sources of evidence make it overwhelmingly
269 likely that the Panubis1.0 assembly provides the correct order and orientation for the
270 sequence. For the remaining large syntenic difference, it is difficult to conclude which
271 one of Panubis1.0 and Panu_3.0 is correct. An example of the nature of this evidence
272 is displayed in Figure 5, which shows that the region starting at ~29.38 Mb and ending
273 at ~44.71 Mb on scaffold NC_018167.2 in Panu_3.0 is inverted relative to the
274 Panubis1.0 assembly. We provide additional information in support of the Panubis1.0
275 assembly from several other regions in Supplementary Figures S1-S5.

276
277

278 **Conclusion**

279

280 The development and commercialization of new technologies by companies such as
281 Illumina, 10x Genomics, Bionano Genomics, Dovetail Genomics and Phase Genomics
282 has enabled researchers to cheaply generate fully *de novo* genome assemblies with
283 high scaffold contiguity (e.g., [40]; [41]; [33]; [36]; [42]). When used in combination with
284 long-read sequences (e.g., from Oxford Nanopore or Pacific Biosciences), these
285 technologies can produce high-quality genome assemblies at a fraction of the cost of
286 traditional clone library based approaches (e.g., [41]; [43]). In this context, our
287 assembly Panubis1.0 provides a 10-fold increase in contig N50 size and a 240-fold
288 increase in scaffold N50 size relative to Panu_3.0 at less than 1% of the reagent cost.
289 The contiguity of this assembly will be especially useful for future studies where
290 knowing the genomic location is important (e.g., hybridization or recombination studies).

291

292 One natural question that arises with any new genome assembly is how one assesses
293 that an assembly is '*correct*'. Indeed, some of the recently published Hi-C based
294 assemblies have not provided any corroborating evidence supporting their assemblies
295 (e.g., [44]). Here, we used three independent sources to provide evidence that 11 out
296 of 12 large syntenic differences are correct in our new baboon assembly (Panubis1.0)
297 relative to the previous assembly Panu_3.0 (Table 3). These include two different
298 sources that contain information about historical patterns of pedigree linkage or linkage
299 disequilibrium across regions. In all, this is substantially more support for our assembly
300 than was produced by previous Hi-C based assemblies (e.g., [41]; [42]; [43]). Finally,

301 we also note that these independent sources of evidence counter any potential criticism
302 of the fact that our genome assembly (using individual '15944' from the SNPRC baboon
303 colony) comes from a different individual from the previous baboon assembly (individual
304 1X1155 from the SNPRC baboon colony). In particular, the linkage and linkage
305 disequilibrium based approaches that we used implicitly average across individuals, and
306 make it much more likely that the differences that we observe are not due to
307 polymorphic structural variation in olive baboons.

308

309 **Availability of supporting data**

310

311 All of the raw sequence data from individual 15944, as well as the Panubis1.0 assembly
312 are available without restriction from NCBI under BioProject PRJNA527874. New RNA-
313 seq data used for genome annotation are available under BioProject PRJNA559725.

314 The genome annotation report and raw files can be found at

315 https://www.ncbi.nlm.nih.gov/genome/annotation_euk/Papio_anubis/104/.

316

317 **Competing interests**

318

319 The authors declare that they have no competing interests.

320

321

322 **Acknowledgements**

323

324 We thank Jenny Tung for providing the RNA sequencing data and some of the high-
325 molecular weight DNA used in this study.

326

327 **Funding**

328

329 The work was supported in part by NIH grants R24 OD017859 (to LAC and JDW), R01
330 GM115433 (to JDW), R01 GM094402 (to YSS), R01 HG005946 (to PYK) and by a
331 Packard Fellowship for Science and Engineering (to YSS). YSS is a Chan Zuckerberg
332 Biohub Investigator.

333

334 **Author contributions**

335

336 JDW, LAC and YSS conceived the project. JG, SD, SS, MLS and PYK generated data
337 for the project. MLS and SSB performed the genome assembly. SSB, MLS, JR, TPV
338 and JDW performed the other analyses. SSB and JDW wrote the manuscript with
339 contributions from all authors.

340

341 **References**

- 342
- 343 1. VandeBerg JL, Williams-Blangero S, Tardif SD. The Baboon in Biomedical Research.
344 Springer Science & Business Media; 2009.
- 345 2. McGill HC, McMahan CA, Kruski AW, Kelley JL, Mott GE. Responses of serum
346 lipoproteins to dietary cholesterol and type of fat in the baboon. *Arteriosclerosis*.
347 1981;1:337–44.
- 348 3. Kushwaha RS, Reardon CA, Lewis DS, Qi Y, Rice KS, Getz GS, et al. Effect of
349 dietary lipids on plasma activity and hepatic mRNA levels of cholesteryl ester transfer
350 protein in high- and low-responding baboons (*Papio* species). *Metabolism*.
351 1994;43:1006–12.
- 352 4. Singh AT, Rainwater DL, Kammerer CM, Sharp RM, Poushesh M, Shelledy WR, et
353 al. Dietary and genetic effects on LDL size measures in baboons. *Arterioscler Thromb*
354 *Vasc Biol*. 1996;16:1448–53.
- 355 5. Kammerer CM, Rainwater DL, Cox LA, Schneider JL, Mahaney MC, Rogers J, et al.
356 Locus controlling LDL cholesterol response to dietary cholesterol is on baboon
357 homologue of human chromosome 6. *Arterioscler Thromb Vasc Biol*. 2002;22:1720–5.
- 358 6. Rainwater DL, Kammerer CM, Mahaney MC, Rogers J, Cox LA, Schneider JL, et al.
359 Localization of genes that control LDL size fractions in baboons. *Atherosclerosis*.
360 2003;168:15–22.
- 361 7. Martin LJ, Blangero J, Rogers J, Mahaney MC, Hixson JE, Carey KD, et al. A
362 quantitative trait locus influencing activin-to-estrogen ratio in pedigreed baboons maps
363 to a region homologous to human chromosome 19. *Hum Biol*. 2001;73:787–800.
- 364 8. Sherwood RJ, Duren DL, Havill LM, Rogers J, Cox LA, Towne B, et al. A
365 genomewide linkage scan for quantitative trait loci influencing the craniofacial complex
366 in baboons (*Papio hamadryas* spp.). *Genetics*. 2008;180:619–28.
- 367 9. Havill LM, Mahaney MC, Cox LA, Morin PA, Joslyn G, Rogers J. A quantitative trait
368 locus for normal variation in forearm bone mineral density in pedigreed baboons maps
369 to the ortholog of human chromosome 11q. *J Clin Endocrinol Metab*. 2005;90:3638–45.
- 370 10. Havill LM, Cox LA, Rogers J, Mahaney MC. Cross-species replication of a serum
371 osteocalcin quantitative trait locus on human chromosome 16q in pedigreed baboons.
372 *Calcif Tissue Int*. 2005;77:205–11.
- 373 11. Rainwater DL, Cox LA, Rogers J, VandeBerg JL, Mahaney MC. Localization of
374 multiple pleiotropic genes for lipoprotein metabolism in baboons. *J Lipid Res*.
375 2009;50:1420–8.

- 376 12. Cox LA, Glenn J, Ascher S, Birnbaum S, VandeBerg JL. Integration of genetic and
377 genomic methods for identification of genes and gene variants encoding QTLs in the
378 nonhuman primate. *Methods*. 2009;49:63–9.
- 379 13. Telenti A, Pierce LCT, Biggs WH, di Iulio J, Wong EHM, Fabani MM, et al. Deep
380 sequencing of 10,000 human genomes. *Proc Natl Acad Sci USA*. 2016;113:11901–6.
- 381 14. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A
382 reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet*.
383 2016;48:1279–83.
- 384 15. Halldorsson BV, Palsson G, Stefansson OA, Jonsson H, Hardarson MT, Eggertsson
385 HP, et al. Characterizing mutagenic effects of recombination through a sequence-level
386 genetic map. *Science*. 2019;363(6425):eaau1043.
- 387 16. Karczewski, Konrad J., et al. Variation across 141,456 human exomes and
388 genomes reveals the spectrum of loss-of-function intolerance across human protein-
389 coding genes. bioRxiv preprint. 2019. <https://doi.org/10.1101/531210>
- 390 17. <https://www.nhlbiwgs.org/>
- 391 18. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial
392 sequencing and analysis of the human genome. *Nature*. 2001;409:860–921.
- 393 19. International Human Genome Sequencing Consortium. Finishing the euchromatic
394 sequence of the human genome. *Nature*. 2004;431:931–45.
- 395 20. International HapMap Consortium. A haplotype map of the human genome. *Nature*.
396 2005;437:1299–320.
- 397 21. International HapMap Consortium. A second generation human haplotype map of
398 over 3.1 million SNPs. *Nature*. 2007 Oct;449(7164):851.
- 399 22. International HapMap 3 Consortium, Altshuler DM, Gibbs RA, Peltonen L, Altshuler
400 DM, Gibbs RA, et al. Integrating common and rare genetic variation in diverse human
401 populations. *Nature*. 2010;467:52–8.
- 402 23. 1000 Genomes Project Consortium, Abecasis GR, Altshuler D, Auton A, Brooks LD,
403 Durbin RM, et al. A map of human genome variation from population-scale sequencing.
404 *Nature*. 2010;467:1061–73.
- 405 24. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo
406 MA, Durbin RM, et al. An integrated map of genetic variation from 1,092 human
407 genomes. *Nature*. 2012;491:56–65.
- 408 25. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP,
409 Kang HM, et al. A global reference for human genetic variation. *Nature*. 2015;526:68–
410 74.

- 411 26. Wall JD, Schlebusch SA, Alberts SC, Cox LA, Snyder-Mackler N, Nevonen KA, et
412 al. Genomewide ancestry and divergence patterns from low-coverage sequencing data
413 reveal a complex history of admixture in wild baboons. *Mol Ecol.* 2016;25:3469–83.
- 414 27. Rogers J, Raveendran M, Harris RA, Mailund T, Leppälä K, Athanasiadis G,
415 Schierup MH, Cheng J, Munch K, Walker JA, Konkel MK. The comparative genomics
416 and complex population history of Papio baboons. *Science Advances.*
417 2019;5(1):eaau6947.
- 418 28. Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. Direct determination of
419 diploid genome sequences. *Genome Res.* 2017;27:757–67.
- 420 29. Ma ZS, Li L, Ye C, Peng M, Zhang YP. Hybrid assembly of ultra-long Nanopore
421 reads augmented with 10x-Genomics contigs: Demonstrated with a human genome.
422 *Genomics.* 2018, in press.
- 423 30. Qin M, Wu S, Li A, Zhao F, Feng H, Ding L, Chang Y, Ruan J. LRScaf: Improving
424 Draft Genomes Using Long Noisy Reads. bioRxiv preprint. 2018.
425 <https://doi.org/10.1101/374868>
- 426 31. Xu GC, Xu TJ, Zhu R, Zhang Y, Li SQ, Wang HW, Li JT. LR_Gapcloser: a tiling
427 path-based gap closer that uses long reads to complete genome assembly.
428 *GigaScience.* 2018;8(1):giy157.
- 429 32. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA,
430 Zeng Q, Wortman J, Young SK, Earl AM. Pilon: an integrated tool for comprehensive
431 microbial variant detection and genome assembly improvement. *PLoS One.*
432 2014;9(11):e112963.
- 433 33. Dudchenko O, Batra SS, Omer AD, Nyquist SK, Hoeger M, Durand NC, et al. De
434 novo assembly of the *Aedes aegypti* genome using Hi-C yields chromosome-length
435 scaffolds. *Science.* 2017;356:92–5.
- 436 34. <https://www.phasegenomics.com/>
- 437 35. Dudchenko O, Shamim MS, Batra S, Durand NC, Musial NT, Mostofa R, Pham M,
438 St Hilaire BG, Yao W, Stamenova E, Hoeger M. The Juicebox Assembly Tools module
439 facilitates de novo assembly of mammalian genomes with chromosome-length scaffolds
440 for under \$1000. bioRxiv preprint. 2018. <https://doi.org/10.1101/254797>
- 441 36. Matthews, Benjamin J., et al. Improved reference genome of *Aedes aegypti* informs
442 arbovirus vector control. *Nature.* 2018;563:501–7.
- 443 37. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO:
444 assessing genome assembly and annotation completeness with single-copy orthologs.
445 *Bioinformatics.* 2015;31:3210–2.

- 446 38. Perelman P, Johnson WE, Roos C, Seuánez HN, Horvath JE, Moreira MAM, et al. A
447 molecular phylogeny of living primates. *PLoS Genet.* 2011;7:e1001342.
- 448 39. Skov L, Schierup MH, Danish Pan Genome Consortium. Analysis of 62 hybrid
449 assembled human Y chromosomes exposes rapid structural changes and high rates of
450 gene conversion. *PLoS Genetics.* 2017;13(8):e1006834.
- 451 40. Mostovoy Y, Levy-Sakin M, Lam J, Lam ET, Hastie AR, Marks P, et al. A hybrid
452 approach for de novo human genome sequence assembly and phasing. *Nat Methods.*
453 2016;13:587–90.
- 454 41. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-
455 molecule sequencing and chromatin conformation capture enable de novo reference
456 assembly of the domestic goat genome. *Nat Genet.* 2017;49:643–50.
- 457 42. Rice, Edward S., et al. Improved genome assembly of American alligator genome
458 reveals conserved architecture of estrogen signaling. *Genome Research.* 2017;27:686–
459 696.
- 460 43. Kalbfleisch, Theodore S., et al. Improved reference genome for the domestic horse
461 increases assembly contiguity and composition. *Communications Biology.* 2018;1:197.
- 462 44. Nuss AB, Sharma A, Gulia-Nuss M. Chicago and Dovetail Hi-C proximity ligation
463 yield chromosome length scaffolds of *Ixodes scapularis* genome. bioRxiv preprint. 2018.
464 <https://doi.org/10.1101/392126>
- 465 45. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore
466 sequencing and assembly of a human genome with ultra-long reads. *Nat Biotechnol.*
467 2018;36:338–45.
- 468 46. Rao, Suhas SP, et al. A 3D map of the human genome at kilobase resolution
469 reveals principles of chromatin looping. *Cell.* 2014;159:1665–1680.
- 470 47. Chan AH, Jenkins PA, Song YS. Genome-wide fine-scale recombination rate
471 variation in *Drosophila melanogaster*. *PLoS Genet.* 2012;8(12):e1003090.
- 472 48. Robinson, Jacqueline A., et al. Analysis of 100 high-coverage genomes from a
473 pedigreed captive baboon colony. *Genome Research.* 2019;29:848–856.
- 474 49. Coop G, Wen X, Ober C, Pritchard JK, Przeworski M. High-resolution mapping of
475 crossovers reveals extensive variation in fine-scale recombination patterns among
476 humans. *Science.* 2008;319:1395–8.

477 **Figure 1. Illustration of our genome assembly strategy.**
478 **Figure 2. The Hi-C map of our Panubis1.0 genome.** Each blue square on the
479 diagonal represents a chromosome-length scaffold. Autosomes are listed first, ordered
480 by size, and the last square corresponds to the X chromosome. The axes are labelled in
481 units of megabases.
482 **Figure 3. Dotplots showing chromosome Y synteny.** A dotplot between rhesus
483 chromosome Y and Panubis1.0 putative chromosome Y is shown on the left, while a
484 dotplot between the chimpanzee chromosome Y and the human chromosome Y is
485 shown on the right. The axes labels are in units of megabases. The phylogenetic
486 distance between baboon and rhesus macaque is similar to that between human and
487 chimpanzee. Hence, the broadly conserved synteny between the rhesus and baboon
488 putative chromosome Y as compared to the synteny between the chimp and human
489 chromosome Y, suggests that the scaffold representing the putative chromosome Y in
490 the Panubis1.0 assembly is indeed capturing at least a large part of chromosome Y.

491 **Figure 4. Dotplots showing alignment of Panu_3.0 reference-assisted**
492 **chromosomes vs. Panubis1.0 chromosome-length scaffolds.** The Panu_3.0
493 assembly is shown on the Y-axis and the Panubis1.0 assembly is shown on the X-axis.
494 Each dot represents the position of a syntenic block between the two assemblies as
495 determined by the nucmer alignment. The color of the dot reflects the orientation of the
496 individual alignments (purple indicates consistent orientation and blue indicates
497 inconsistent orientation). The dotplots illustrate that there are chromosomes containing
498 large inversions and translocations in the Panu_3.0 assembly with respect to the
499 Panubis1.0 assembly.

500 **Figure 5. Evidence for misassembly on chromosome NC_018167.2 in Panu_3.0.**

501 **a)** Bionano optical map alignment to the Panu_3.0 assembly demonstrates there is an
502 inversion on chromosome NC_018167.2 beginning at ~29.38 Mb and ending at ~44.71
503 Mb. **b)** Estimates of the population recombination rate ρ near the potential syntenic
504 breaks of the inversion identified on chromosome NC_018167.2. **c)** Shown on the x-
505 axis is positions along chromosome NC_018167.2 in Panu_3.0 where each row
506 represents one of the 9 offsprings of sire 10173. Switches between red and blue within
507 a row represent a recombination event. The two vertical black lines represent locations
508 where three or more recombinations occur at the same locus indicating a potential
509 misassembly.

510 **Figure 6. Pedigree of baboons used in linkage analysis.**

Assembly	10x	10x contigs	10x contigs + Nanopore scaffolding	10x contigs + Nanopore scaffolding + Nanopore gap filling	10x contigs + Nanopore scaffolding + Nanopore gap filling + Illumina polishing	Panubis1.0
Total Length of Scaffolds	2,894,992,835	2,809,352,255	2,871,292,557	2,871,210,925	2,870,847,162	2,871,135,062
Number of Scaffolds	24,429	87,632	15,803	15,803	15,803	12,976
Scaffold N50	20,460,278	84,258	1,695,573	1,695,772	1,695,642	140,274,886
Total Gap Length	85,640,580	0	50,344,034	2,030,908	2,030,908	2,318,808
Total Length of Contigs	2,809,352,255	2,809,352,255	2,820,948,523	2,869,180,017	2,868,816,254	2,868,816,254
Number of Contigs	87,632	87,632	62,252	17,004	17,004	17,055
Contig N50	84,258	84,258	134,222	1,469,760	1,469,602	1,461,245
BUSCO Score	92.70%	74.20%	92.70%	92.90%	93.00%	93.00%

511
512 **Table 1. Assembly statistics for each step of the adopted assembly strategy.**
513 Total Length of Scaffolds is the sum of lengths of scaffolds (including A, C, G, T and N)
514 in each scaffold. Total Gap Length is the total number of N's in the assembly.
515 Total Length of Contigs is the sum of the number of sequenced base pairs (including
516 only A, C, G and T) in each scaffold. BUSCO provides a way of measuring the
517 presence of genes conserved in mammals [37]. Since BUSCO reports complete genes
518 and fragmented genes, the BUSCO Score is the fraction of complete mammalian
519 BUSCO genes found in the assembly.

520

Assembly	Panubis1.0	Panu_3.0
Total Length of Scaffolds	2,871,135,062	2,959,373,024
Number of Scaffolds	12,976	63,235
Scaffold N50	140,274,886	585,721
Total Gap Length	2,318,808	22,434,732
Total Length of Contigs	2,868,816,254	2,937,001,527
Number of Contigs	17,055	122,216
Contig N50	1,461,245	138,819
BUSCO Score	93.00%	93.40%

521

522

523

Table 2. Comparison of Panubis1.0 with Panu_3.0 assemblies.

524

Panu_3.0 chromosome	Panu_3.0 Start	Panu_3.0 End	Panu_2.0 Start	Panu_2.0 End	Type	Linkage support	BNG support	LDhelmet support
NC_018164.2	88.05	104.99	87.61	104.98	Inv	start ¹	yes	unknown ¹
NC_018167.2	29.38	44.71	29.25	44.53	Inv	start + end	yes	start + end
NC_018156.2	4.04	8.67	4.18	8.63	Inv	no	yes ²	no
NC_018162.2	82.42	86.47	81.91	84.01	Trans	start + end	no ³	no
NC_018166.2	104.28	108.05	103.66	107.44	Inv	no	yes	no
NC_018165.2	15.93	19.48	15.85	19.40	Inv	no	no	no
NC_018166.2	96.94	100.12	96.39	99.54	Trans	start + end	yes ⁴	start + end
NC_018160.2	36.05	36.75	35.88	36.55	Trans	no	yes ⁴	start
NC_018163.2	23.19	23.66	0	0.47	Trans	no	yes ²	no
NC_018164.2	4.05	4.49	3.99	4.45	Trans	no ⁵	yes	no
NC_018165.2	100.91	101.18	100.31	100.59	Trans	no	yes	no
NC_018152.2	166.73	166.89	169.86	170.10	Trans	start + end	yes	end

525

526

527

Table 3. Potential large (>100 Kb) assembly errors in Panu_3.0, ordered by size.

Note that a 'no' in the 'Linkage support' or 'LDhelmet support' columns is inconclusive, and should not be interpreted as support for the Panu_3.0 assembly being correct.

531

¹ Unable to determine whether linkage and LDhelmet provide support at the end breakpoint due to a lack of synteny between Panu_2.0 and Panu_3.0

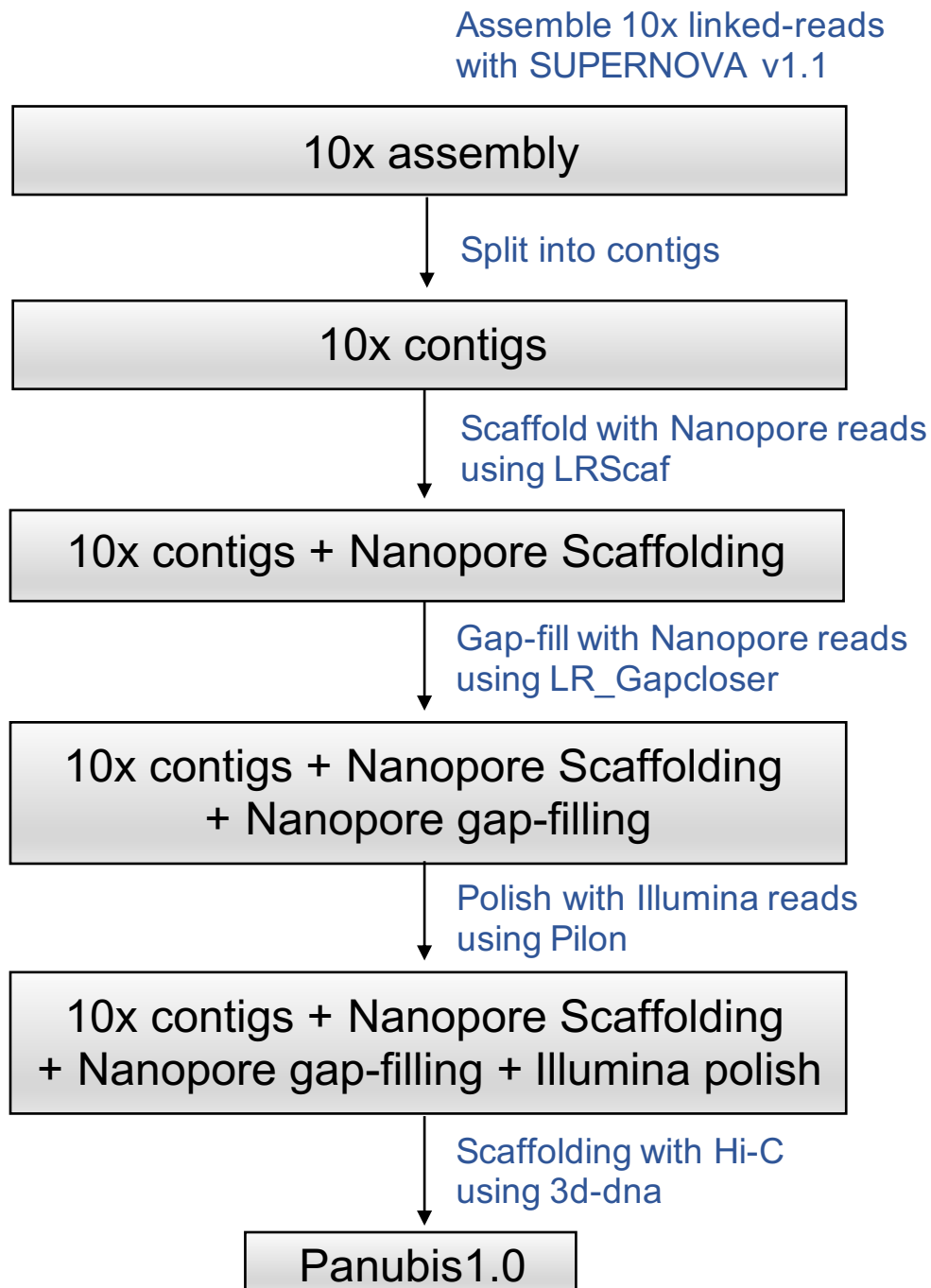
² Panu_2.0 assembly appears to be correct

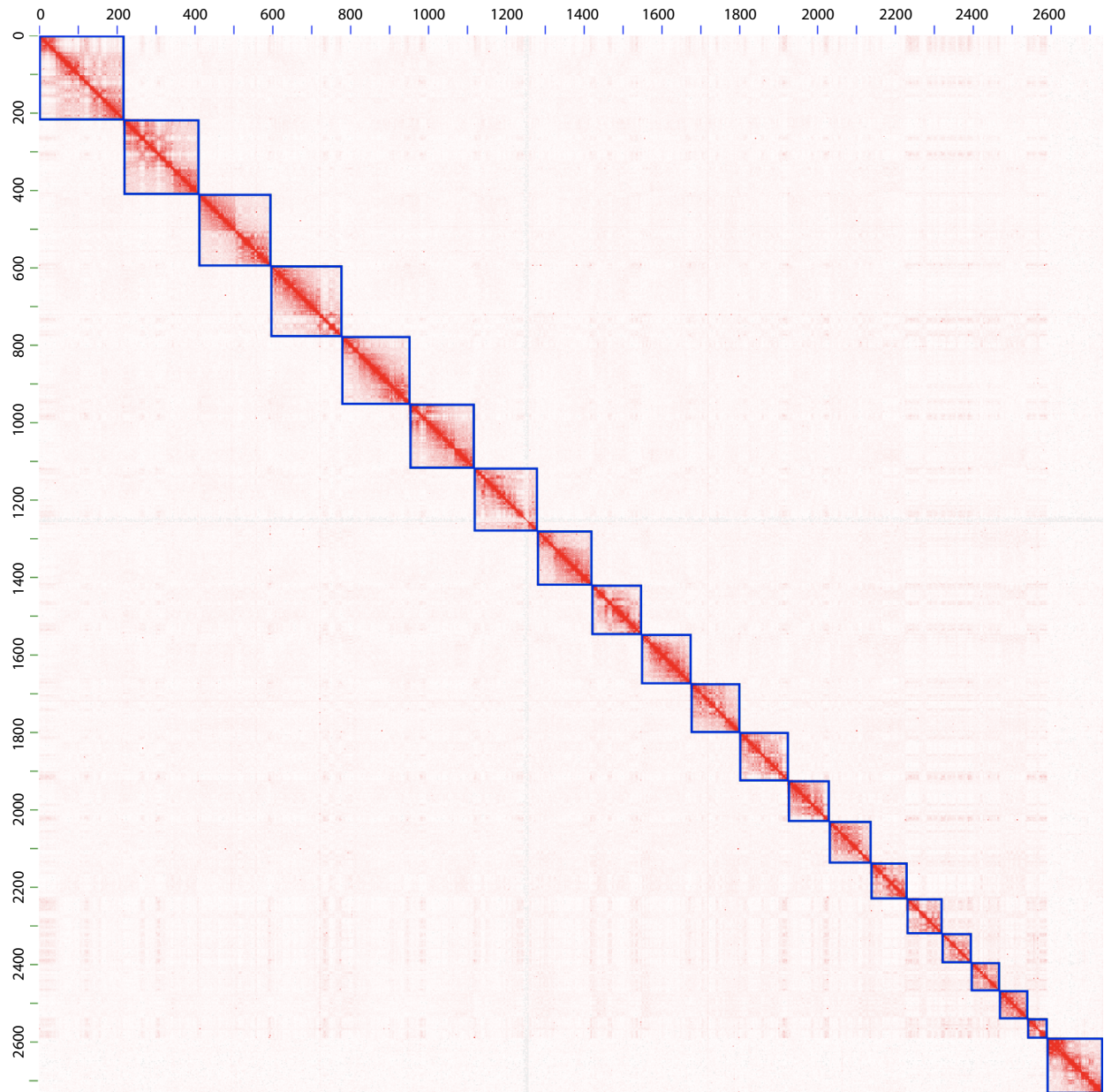
³ BNG maps do not support a translocation with these breakpoints. However, they do support a potential large SV at the starting breakpoint

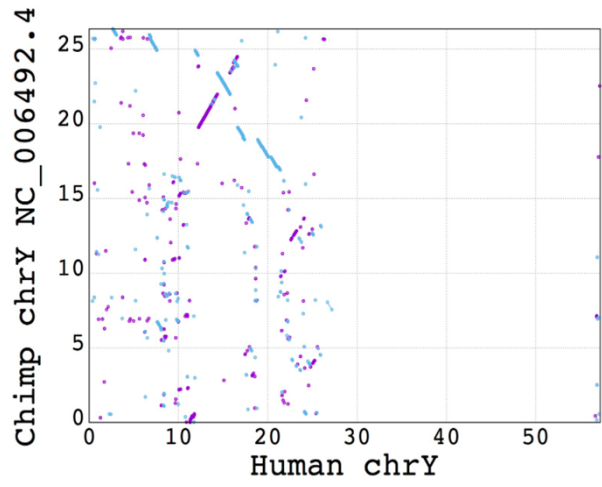
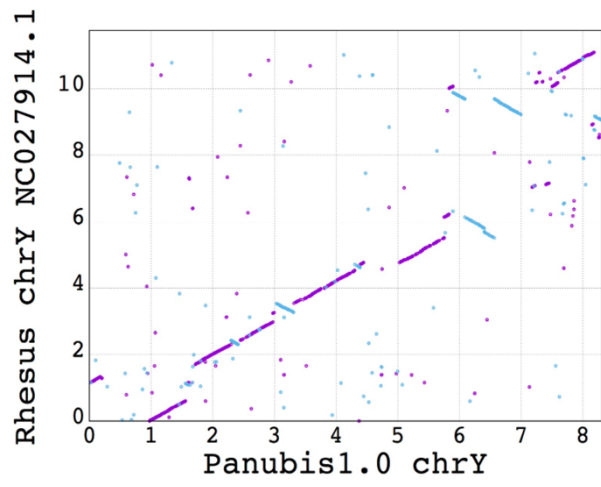
⁴ BNG maps support the presence of a large SV, which may be a translocation

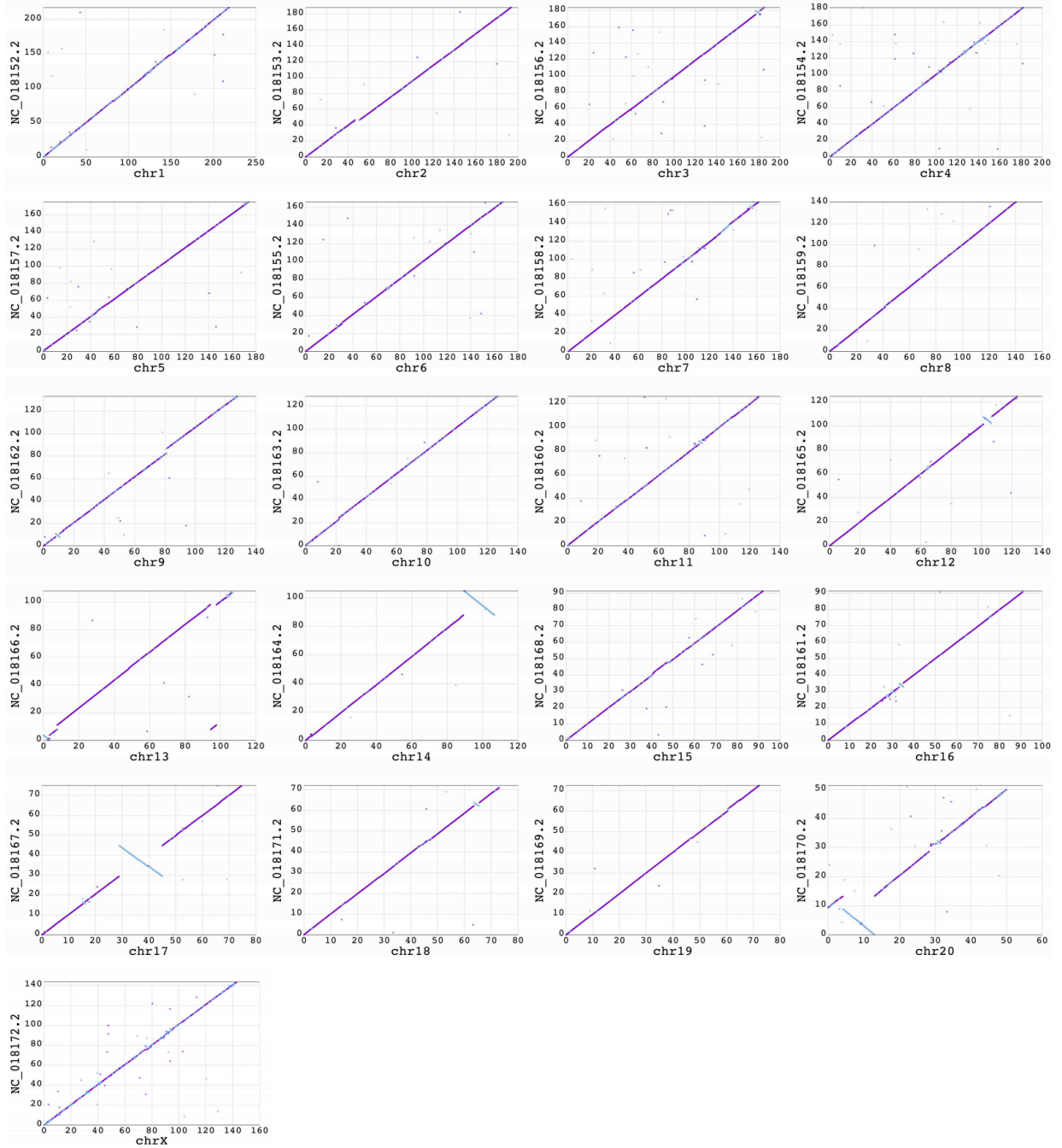
⁵ Linkage data suggests a potential polymorphic inversion (in 16413) partially overlapping with this interval

539

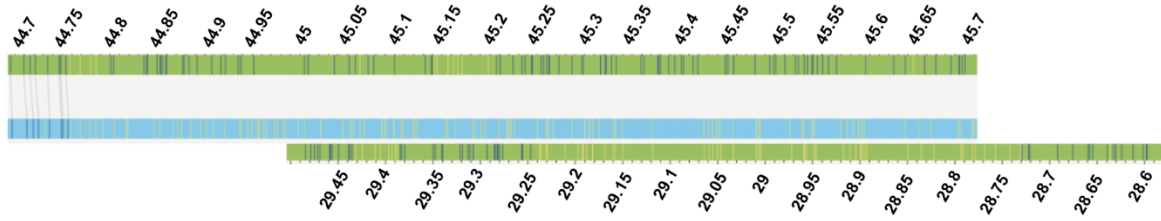




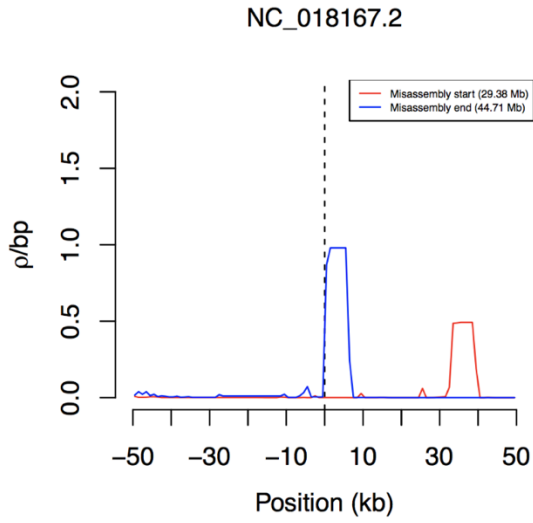




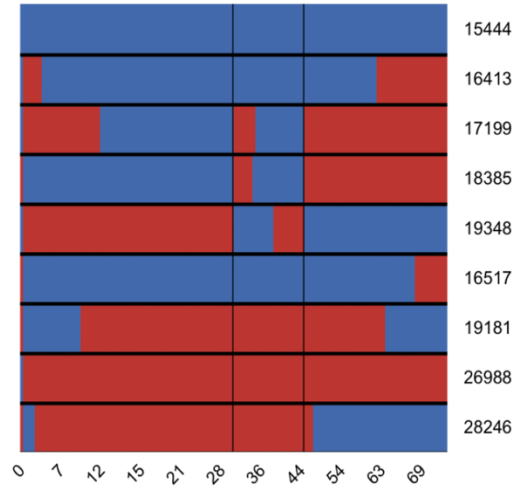
a NC_018167.2: 29.38-44.71 Inversion

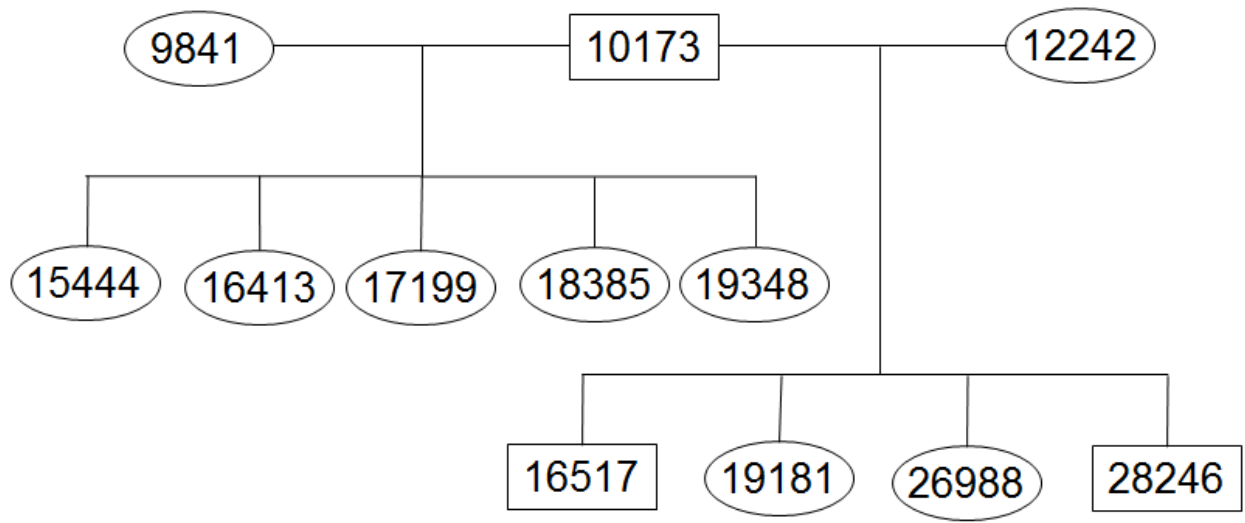


b



c







Click here to access/download
Supplementary Material
Supplementary_Material.pdf

